

PCT/JP 00/00423  
27.01.00

日 本 国 特 許 庁

PATENT OFFICE  
JAPANESE GOVERNMENT

EU

REC'D 14 FEB 2000

WIPO PCT

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日

Date of Application: 1999年 1月29日

出 願 番 号

Application Number: 平成11年特許願第023064号

出 願 人

Applicant (s): ソニー株式会社

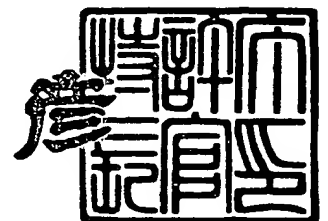
PRIORITY  
DOCUMENT

SUBMITTED OR TRANSMITTED IN  
COMPLIANCE WITH RULE 17.1(a) OR (b)

1999年11月19日

特許庁長官  
Commissioner,  
Patent Office

近 藤 隆 彦



出証番号 出証特平11-3080174

【書類名】 特許願  
 【整理番号】 9801128002  
 【あて先】 特許庁長官殿  
 【国際特許分類】 G01B 11/24  
 【発明者】

【住所又は居所】 東京都品川区北品川 6 丁目 7 番 3 5 号 ソニー株式会社  
 内

【氏名】 トビー ウォーカー  
 【発明者】

【住所又は居所】 東京都品川区北品川 6 丁目 7 番 3 5 号 ソニー株式会社  
 内

【氏名】 松原 弘  
 【特許出願人】  
 【識別番号】 000002185  
 【氏名又は名称】 ソニー株式会社  
 【代表者】 出井 伸之

【代理人】  
 【識別番号】 100067736  
 【弁理士】

【氏名又は名称】 小池 晃  
 【選任した代理人】

【識別番号】 100086335  
 【弁理士】  
 【氏名又は名称】 田村 榮一

【選任した代理人】  
 【識別番号】 100096677  
 【弁理士】  
 【氏名又は名称】 伊賀 誠司

【手数料の表示】

【予納台帳番号】 019530

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

---

【包括委任状番号】 9707387

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 信号処理方法及び映像音声処理装置

【特許請求の範囲】

【請求項1】 供給された信号の内容の意味構造を反映するパターンを検出して解析する信号処理方法であって、

上記信号を構成する連続したフレームのひと続きから形成されるセグメントから、その特性を表す少なくとも1つ以上の特徴情報を抽出する特徴抽出工程と、

上記特徴情報を用いて、上記特徴情報のそれぞれ毎に、上記セグメントの対の間の類似性を測定する測定基準を算出して、この測定基準により上記セグメントの対の間の類似性を測定する類似性測定工程と、

上記特徴情報と上記測定基準とを用いて、上記セグメントのうち、互いの時間的距離が所定の時間閾値以内であるとともに、互いの類似性が所定の非類似性閾値以上である2つのセグメントを検出し、上記信号の内容の意味構造を反映し、時間的に連続するセグメントからなるシーンにまとめるグループ化工程とを備えること

を特徴とする信号処理方法。

【請求項2】 上記信号とは、ビデオデータにおける映像信号と音声信号との少なくとも1つであること

を特徴とする請求項1記載の信号処理方法。

【請求項3】 上記特徴抽出工程により、複数個に分割された単一のセグメント内の異なる時点における上記特徴情報の単一の統計的な代表値を選択して抽出すること

を特徴とする請求項1記載の信号処理方法。

【請求項4】 上記非類似性閾値は、複数個のセグメントの対の間の類似性の統計値を用いて決定されること

を特徴とする請求項1記載の信号処理方法。

【請求項5】 上記セグメントのうち、上記グループ化工程にてシーンにまとめられなかった少なくとも1つ以上のセグメントを、単一のシーンとしてまとめること

を特徴とする請求項 1 記載の信号処理方法。

【請求項 6】 上記グループ化工程により得られた任意の特徴情報に関するシーンと、上記グループ化工程により得られた上記任意の特徴情報とは異なる特徴情報に関する少なくとも 1 つ以上のシーンとを単一にまとめること

を特徴とする請求項 1 記載の信号処理方法。

【請求項 7】 上記グループ化工程により得られた上記映像信号における少なくとも 1 つ以上のシーンと、上記グループ化工程により得られた上記音声信号における少なくとも 1 つ以上のシーンとを単一にまとめること

を特徴とする請求項 2 記載の信号処理方法。

【請求項 8】 供給されたビデオ信号の内容の意味構造を反映する映像及び／又は音声のパターンを検出して解析する映像音声処理装置であって、

上記ビデオ信号を構成する連続した映像及び／又は音声フレームのひと続きから形成される映像及び／又は音声セグメントから、その特性を表す少なくとも 1 つ以上の特徴情報を抽出する特徴抽出手段と、

上記特徴情報を用いて、上記特徴情報のそれぞれ毎に、上記映像及び／又は音声セグメントの対の間の類似性を測定する測定基準を算出して、この測定基準により上記映像及び／又は音声セグメントの対の間の類似性を測定する類似性測定手段と、

上記特徴情報と上記測定基準とを用いて、上記映像及び／又は音声セグメントのうち、互いの時間的距離が所定の時間閾値以内であるとともに、互いの類似性が所定の非類似性閾値以上である 2 つの映像及び／又は音声セグメントを検出し、上記ビデオ信号の内容の意味構造を反映し、時間的に連続する映像及び／又は音声セグメントからなるシーンにまとめるグループ化手段とを備えること

を特徴とする映像音声処理装置。

【請求項 9】 上記特徴抽出手段は、複数個に分割された単一の映像及び／又は音声セグメント内の異なる時点における上記特徴情報の単一の統計的な代表値を選択して抽出すること

を特徴とする請求項 8 記載の映像音声処理装置。

【請求項 10】 上記非類似性閾値は、複数の映像及び／又は音声セグメントの対の間の類似性の統計値を用いて決定されること

を特徴とする請求項 8 記載の映像音声処理装置。

【請求項 11】 上記映像及び／又は音声セグメントのうち、上記グループ化手段によりシーンにまとめられなかった少なくとも 1 つ以上の映像及び／又は音声セグメントを、単一のシーンとしてまとめること

を特徴とする請求項 8 記載の映像音声処理装置。

【請求項 12】 上記グループ化手段により得られた任意の特徴情報に関するシーンと、上記グループ化手段により得られた上記任意の特徴情報とは異なる特徴情報に関する少なくとも 1 つ以上のシーンとを単一にまとめること

を特徴とする請求項 8 記載の映像音声処理装置。

【請求項 13】 上記グループ化工程により得られた上記ビデオ信号の映像信号における少なくとも 1 つ以上のシーンと、上記グループ化工程により得られた上記ビデオ信号の音声信号における少なくとも 1 つ以上のシーンとを単一にまとめること

を特徴とする請求項 8 記載の映像音声処理装置。

#### 【発明の詳細な説明】

##### 【0001】

#### 【発明の属する技術分野】

本発明は、信号の基礎となる意味構造を反映するパターンを検出して解析する信号処理方法及びビデオ信号の基礎となる意味構造を反映する映像及び／又は音声のパターンを検出して解析する映像音声処理装置に関する。

##### 【0002】

#### 【従来の技術】

例えばビデオデータに録画されたテレビ番組といった大量の異なる映像データにより構成される映像アプリケーションの中から、興味のある部分等の所望の部分を探して再生したい場合がある。

##### 【0003】

このように、所望の映像内容を抽出するための一般的な技術としては、アプリ

ケーションの主要場面を描いた一連の映像を並べて作成されたパネルであるストーリーボードがある。このストーリーボードは、ビデオデータをいわゆるショットに分解し、各ショットにおいて代表される映像を表示したものである。このような映像抽出技術は、そのほとんどが、例えば“G. Ahanger and T.D.C. Little, A survey of technologies for parsing and indexing digital video, J. of Visual Communication and Image Representation 7:28-4, 1996”に記載されているように、ビデオストラクチャからショットを自動的に検出して抽出するものである。

#### 【0004】

##### 【発明が解決しようとする課題】

ところで、例えば代表的な30分のテレビ番組中には、数百ものショットが含まれている。そのため、上述した従来の映像抽出技術においては、ユーザが抽出された膨大な数のショットを並べたストーリーボードを調べる必要があり、このようなストーリーボードを理解するにはユーザに大きな負担を強いる必要があった。また、従来の映像抽出技術においては、例えば話し手の変化に応じて交互に2者を撮影した会話場面におけるショットは、冗長のものが多いという問題があった。このように、ショットは、ビデオストラクチャから抽出する対象としては階層が低すぎて必要のない情報量が多く、このようなショットを抽出する従来の映像抽出技術は、ユーザにとって利便性のよいものとはいえなかった。

#### 【0005】

また、他の映像抽出技術としては、例えば“A. Merlino, D. Morey and M. Maybury, Broadcast news navigation using story segmentation, Proc. of ACM Multimedia 97, 1997”や特開平10-136297号公報に記載されているように、ニュースやフットボールゲームといった特定の内容ジャンルに関する非常に専門的な知識を用いるものがある。しかしながら、この従来の映像抽出技術は、目的のジャンルに関しては良好な結果を得ることができるが、他の種類のジャンルには全く役に立つものではなく、ジャンルに限定されて容易に一般化することができないという問題があった。

## 【0006】

さらに、他の映像抽出技術としては、例えばU.S. Patent #5,708,767号公報に記載されているように、いわゆるストーリーユニットを抽出するものがある。しかしながら、この従来の映像抽出技術は、完全に自動化されたものではなく、どのショットが同じ内容を示すものであるかを決定するために、ユーザの介入が必要であった。また、この従来の映像抽出技術は、処理に要する計算が複雑であるとともに、適用対象として映像情報のみに限定されるといった問題もあった。

## 【0007】

さらにまた、他の映像抽出技術としては、例えば特開平9-214879号公報に記載されているように、ショット検出と無音部分検出とを組み合わせることによりショットを識別するものがある。しかしながら、この従来の映像抽出技術は、2つの無音ポイント間にある1組のショットを識別するものであって、無音ポイントがショットの境界に対応する場合のみに限定されたものであった。

## 【0008】

また、他の映像抽出技術としては、例えば“H. Aoki, S. Shimotsuji and O. Hori, A shot classification method to select effective key-frames for video browsing, IPSJ Human Interface SIG Notes, 7:43-50, 1996”や特開平9-93588号公報に記載されているように、ストーリーボードにおける表示の冗長を低減するために反復されたショットを検出するものがある。しかしながら、この従来の映像抽出技術は、映像情報のみに適用できるものであり、音声情報に適用できるものではなかった。

## 【0009】

本発明は、このような実情に鑑みてなされたものであり、上述した従来の映像抽出技術の問題を解決し、種々のビデオデータにおける高レベルのビデオストラクチャを抽出する信号処理方法及び映像音声処理装置を提供することを目的とするものである。

## 【0010】

## 【課題を解決するための手段】

上述した目的を達成する本発明にかかる信号処理方法は、供給された信号の内



容の意味構造を反映するパターンを検出して解析する信号処理方法であって、信号を構成する連続したフレームのひと続きから形成されるセグメントから、その特性を表す少なくとも1つ以上の特徴情報を抽出する特徴抽出工程と、特徴情報を用いて、特徴情報のそれぞれ毎に、セグメントの対の間の類似性を測定する測定基準を算出して、この測定基準によりセグメントの対の間の類似性を測定する類似性測定工程と、特徴情報と測定基準とを用いて、セグメントのうち、互いの時間的距離が所定の時間閾値以内であるとともに、互いの類似性が所定の非類似性閾値以上である2つのセグメントを検出し、信号の内容の意味構造を反映し、時間的に連続するセグメントからなるシーンにまとめるグループ化工程とを備えることを特徴としている。

## 【0011】

このような本発明にかかる信号処理方法は、信号において類似したセグメントを検出してシーンにまとめる。

## 【0012】

また、上述した目的を達成する本発明にかかる映像音声処理装置は、供給されたビデオ信号の内容の意味構造を反映する映像及び／又は音声のパターンを検出して解析する映像音声処理装置であって、ビデオ信号を構成する連続した映像及び／又は音声フレームのひと続きから形成される映像及び／又は音声セグメントから、その特性を表す少なくとも1つ以上の特徴情報を抽出する特徴抽出手段と、特徴情報を用いて、特徴情報のそれぞれ毎に、映像及び／又は音声セグメントの対の間の類似性を測定する測定基準を算出して、この測定基準により映像及び／又は音声セグメントの対の間の類似性を測定する類似性測定手段と、特徴情報と測定基準とを用いて、映像及び／又は音声セグメントのうち、互いの時間的距離が所定の時間閾値以内であるとともに、互いの類似性が所定の非類似性閾値以上である2つの映像及び／又は音声セグメントを検出し、ビデオ信号の内容の意味構造を反映し、時間的に連続する映像及び／又は音声セグメントからなるシーンにまとめるグループ化手段とを備えることを特徴としている。

## 【0013】

このように構成された本発明にかかる映像音声処理装置は、ビデオ信号におい

て類似した映像及び／又は音声セグメントを検出してまとめ、シーンとして出力する。

【0014】

【発明の実施の形態】

以下、本発明を適用した具体的な実施の形態について図面を参照しながら詳細に説明する。

【0015】

本発明を適用した実施の形態は、録画されたビデオデータから所望の内容を自動的に探し出して抽出する映像音声処理装置である。この映像音声処理装置の具体的な説明を行う前に、ここではまず本発明において対象とするビデオデータに関する説明を行う。

【0016】

本発明において対象とするビデオデータについては、図1に示すようにモデル化し、フレーム、セグメント、シーンの3つのレベルに階層化されたストラクチャを有するものとする。すなわち、ビデオデータは、最下位層において、一連のフレームにより構成される。また、ビデオデータは、フレームの1つ上の階層として、連続するフレームのひと続きから形成されるセグメントにより構成される。さらに、ビデオデータは、最上位層において、このセグメントを意味のある関連に基づきまとめて形成されるシーンにより構成される。

【0017】

このビデオデータは、映像及び音声の両方の情報を含むものである。すなわち、このビデオデータにおいてフレームは、単一の静止画像である映像フレームと、一般に数十～数百ミリセカンド／長といった短時間においてサンプルされた音声情報を表す音声フレームとを含むものである。

【0018】

また、セグメントは、単一のカメラにより連続的に撮影された映像フレームのひと続きからなり一般にショットと呼ばれる映像セグメントと音声セグメントとを含むものである。このセグメントは、ビデオストラクチャにおける基本となる単位である。セグメントのうち、音声セグメントは、多くの定義が可能とされ、

例えば、以下に示すようなものが考えられる。例えば、音声セグメントは、一般によく知られている方法により検出されたビデオデータ中の無音期間により境界を定められて形成されるものがある。また、音声セグメントは、“D. Kimber and L. Wilcox, Acoustic Segmentation for Audio Browsers, Xerox Parc Technical Report”に記載されているように、例えば、音声、音楽、ボイド、無音等のように少数のカテゴリに分類された音声フレームのひと続きから形成されるものがある。さらに、音声セグメントは、“S. Pfeiffer, S. Fischer and E. Wolfgang, Automatic Audio Content Analysis, Proceeding of ACM Multimedia 96, Nov. 1996, pp21-30”に記載されているように、2枚の連続する音声フレーム間のある特性における大きな変化を検出する音声カット検出を用いて決定されるものがある。

#### 【0019】

さらに、シーンは、ビデオデータの内容を意味のあるより高いレベルでまとめてグループ化するために、映像セグメント（ショット）検出或いは音声セグメント検出により映像フレーム或いは音声フレームを捕捉したセグメントを、例えばセグメント内の知覚的アクティビティ量といったセグメントの特性を表す特徴情報であるフィーチャを用いて意味のあるまとまりにグループ化したものである。シーンは、主観的でありビデオデータの内容のジャンルに依存するものであるが、ここでは、映像又は音声における幾つかの関連性のあるフィーチャが互いに類似性を示す映像セグメント又は音声セグメントの反復パターンをグループ化したものとする。具体的には、図2に示すように、2人の話し手が互いに会話している場面で、映像セグメントは、話し手に応じて交互に現れる。このような反復パターンを有するビデオデータにおいて、一方の話し手における一連の映像セグメントAと、他方の話し手における一連の映像セグメントBとは、グループにまとめられて1つのシーンを構成する。このような反復パターンは、ビデオデータにおける高いレベルでの意味のあるストラクチャと非常に関係があり、シーンは、このようなビデオデータにおける高いレベルでの意味のあるまとまりを示すものである。

## 【0020】

本発明を適用した実施の形態として図3に示す映像音声処理装置10は、上述したビデオデータにおけるセグメントのフィーチャを用いてセグメント間の類似性を測定し、これらのセグメントをシーンにまとめてビデオストラクチャを自動的に抽出するものであり、映像セグメント及び音声セグメントの両方に適用できるものである。

## 【0021】

映像音声処理装置10は、同図に示すように、入力したビデオデータのストリームを映像、音声又はこれらの両方のセグメントに分割するビデオ分割部11と、ビデオデータの分割情報を記憶するビデオセグメントメモリ12と、各映像セグメントにおけるフィーチャを抽出する特徴抽出手段である映像フィーチャ抽出部13と、各音声セグメントにおけるフィーチャを抽出する特徴抽出手段である音声フィーチャ抽出部14と、映像セグメント及び音声セグメントのフィーチャを記憶するセグメントフィーチャメモリ15と、映像セグメント及び音声セグメントをシーンにまとめるグループ化手段であるシーン検出部16と、2つのセグメント間の類似性を測定する類似性測定手段であるフィーチャ類似性測定部17とを備える。

## 【0022】

ビデオ分割部11は、例えば、MPEG1 (Moving Picture Experts Group phase 1) やMPEG2 (Moving Picture Experts Group phase 2)、いわゆるDV (Digital Video) のような圧縮ビデオデータフォーマットを含む種々のデジタル化されたフォーマットにおける映像データと音声データとからなるビデオデータのストリームを入力し、このビデオデータを映像、音声又はこれらの両方のセグメントに分割するものである。このビデオ分割部11は、入力したビデオデータが圧縮フォーマットであった場合、この圧縮ビデオデータを完全伸張することなく直接処理することができる。ビデオ分割部11は、入力したビデオデータを処理し、映像セグメントと音声セグメントとに分割する。また、ビデオ分割部11は、入力したビデオデータを分割した結果である分割情報を後段のビデオセグメントメモリ12に出力する。さらに、ビデオ分割部11は、映像セグメン

トと音声セグメントとに応じて、分割情報を後段の映像フィーチャ抽出部 13 及び音声フィーチャ抽出部 14 に出力する。

【0023】

ビデオセグメントメモリ 12 は、ビデオ分割部 11 から供給されたビデオデータの分割情報を記憶する。また、ビデオセグメントメモリ 12 は、後述するシーン検出部 16 からの問い合わせに応じて、分割情報をシーン検出部 16 に出力する。

【0024】

映像フィーチャ抽出部 13 は、ビデオ分割部 11 によりビデオデータを分割して得た各映像セグメント毎のフィーチャを抽出する。映像フィーチャ抽出部 13 は、圧縮映像データを完全伸張することなく直接処理することができる。映像フィーチャ抽出部 13 は、抽出した各映像セグメントのフィーチャを後段のセグメントフィーチャメモリ 15 に出力する。

【0025】

音声フィーチャ抽出部 14 は、ビデオ分割部 11 によりビデオデータを分割して得た各音声セグメント毎のフィーチャを抽出する。音声フィーチャ抽出部 14 は、圧縮音声データを完全伸張することなく直接処理することができる。音声フィーチャ抽出部 14 は、抽出した各音声セグメントのフィーチャを後段のセグメントフィーチャメモリ 15 に出力する。

【0026】

セグメントフィーチャメモリ 15 は、映像フィーチャ抽出部 13 及び音声フィーチャ抽出部 14 からそれぞれ供給された映像セグメント及び音声セグメントのフィーチャを記憶する。セグメントフィーチャメモリ 15 は、後述するフィーチャ類似性測定部 17 からの問い合わせに応じて、記憶しているフィーチャの値やセグメントをフィーチャ類似性測定部 17 に出力する。

【0027】

シーン検出部 16 は、ビデオセグメントメモリ 12 に保持された分割情報と、1 対のセグメント間の類似性とを用いて映像セグメント及び音声セグメントをそれぞれシーンにまとめる。シーン検出部 16 は、グループ内の各セグメントから

開始して、セグメント群の中から類似しているセグメントの反復パターンを検出し、このようなセグメントを同一シーンとしてまとめてグループ化する。このシーン検出部 16 は、あるシーンにおけるセグメントをまとめてグループを徐々に大きくしていき、全てのセグメントをグループ化するまで処理を行い、最終的に検出シーンを生成して出力する。シーン検出部 16 は、フィーチャ類似性測定部 17 を用いて、2つのセグメントがどの程度類似しているかを判断する。

#### 【0028】

フィーチャ類似性測定部 17 は、2つのセグメント間の類似性を測定する。フィーチャ類似性測定部 17 は、あるセグメントに関するフィーチャの値を検索するようにセグメントフィーチャメモリ 15 に問いかける。

#### 【0029】

時間的に近接して反復している類似したセグメントは、ほぼ同一シーンの一部であるため、映像音声処理装置 10 は、このようなセグメントを検出してグループ化していくことによって、シーンを検出する。このような映像音声処理装置 10 は、図 4 に概略を示すような一連の処理を行うことによって、シーンを検出する。

#### 【0030】

まず、映像音声処理装置 10 は、同図に示すように、ステップ S1 において、ビデオ分割を行う。すなわち、映像音声処理装置 10 は、ビデオ分割部 11 に入力されたビデオデータを映像セグメント又は音声セグメントのいずれか、或いは可能であればその両方に分割する。映像音声処理装置 10 は、適用するビデオ分割方法に特に前提要件を設けない。例えば、映像音声処理装置 10 は、“G. Ahanger and T.D.C. Little, A survey of technologies for parsing and indexing digital video, J. of Visual Communication and Image Representation 7:28-4, 1996” に記載されているような方法によりビデオ分割を行う。このようなビデオ分割の方法は、当該技術分野ではよく知られたものであり、映像音声処理装置 10 は、いかなるビデオ分割方法も適用できるものとする。

#### 【0031】

次に、映像音声処理装置 10 は、ステップ S2 において、フィーチャの抽出を

行う。すなわち、映像音声処理装置 10 は、映像フィーチャ抽出部 13 や音声フィーチャ抽出部 14 により各セグメントについて 1 組のフィーチャを計算し、セグメントの特性を表す。映像音声処理装置 10 は、例えば、各セグメントの時間長や、カラーヒストグラムやテクスチャフィーチャといった映像フィーチャや、周波数解析結果、レベル、ピッチといった音声フィーチャやアクティビティ測定結果等を適用可能なフィーチャとして計算する。勿論、映像音声処理装置 10 は、適用可能なフィーチャとしてこれらに限定されるものではない。

#### 【0032】

さらに、映像音声処理装置 10 は、ステップ S3 において、フィーチャを用いたセグメントの類似性測定を行う。すなわち、映像音声処理装置 10 は、フィーチャ類似性測定部 17 により非類似性測定を行い、その測定基準により 2 つのセグメントがどの程度類似しているかを測定する。映像音声処理装置 10 は、先のステップ S2 において抽出したフィーチャを用いて、非類似性測定基準を計算する。

#### 【0033】

そして、映像音声処理装置 10 は、ステップ S4 において、セグメントのグループ化を行う。すなわち、映像音声処理装置 10 は、先のステップ S3 において計算した非類似性測定基準と、先のステップ S2 において抽出したフィーチャとを用いて、時間的に近接して類似したセグメントを繰り返しまとめ、これらのセグメントを重ねてグループに合わせる。映像音声処理装置 10 は、このようにして最終的に生成されたグループを検出シーンとして出力する。

#### 【0034】

このような一連の処理を経ることによって、映像音声処理装置 10 は、ビデオデータからシーンを検出することができる。したがって、ユーザは、この結果を用いることによって、ビデオデータの内容を要約したり、ビデオデータ中の興味のあるポイントに迅速にアクセスしたりすることが可能となる。

#### 【0035】

以下、同図に示した映像音声処理装置 10 における処理を各工程毎により詳細に説明していく。

## 【0036】

まず、ステップS1におけるビデオ分割について説明する。映像音声処理装置10は、ビデオ分割部11に入力されたビデオデータを映像セグメント又は音声セグメントのいずれか、或いは可能であればその両方に分割するが、このビデオデータにおけるセグメントの境界を自動的に検出するための技術は、多くのものがあり、映像音声処理装置10において、このビデオ分割方法に特別な前提要件を設けないことは上述した通りである。しかしながら、映像音声処理装置10において、後の工程によるシーン検出の精度は、本質的に、基礎となるビデオ分割の精度に依存する。なお、映像音声処理装置10におけるシーン検出は、ある程度ビデオ分割時のエラーを許容することができる。特に、映像音声処理装置10において、ビデオ分割は、セグメント検出が不十分である部分よりも、セグメント検出が過度である部分においてエラーを生じる方が好ましい。映像音声処理装置10は、類似したセグメントの検出が過度である結果である限り、たとえセグメント検出が過度でなかったとしても、一般に、シーン検出の際に検出過度であるセグメントを同一シーンとしてまとめることができる。

## 【0037】

つぎに、ステップS2におけるフィーチャ抽出について説明する。フィーチャは、セグメントの内容を表すとともに、異なるセグメント間の類似性を測定するためのデータを供給するセグメントの属性である。映像音声処理装置10は、映像フィーチャ抽出部13や音声フィーチャ抽出部14により各セグメントについて1組のフィーチャを計算し、セグメントの特性を表す。映像音声処理装置10は、いかなるフィーチャの具体的詳細にも依存するものではない。映像音声処理装置10において用いて効果的であると考えられるフィーチャとしては、例えば以下に示す映像フィーチャ、音声フィーチャ、映像音声共通フィーチャのようなものがある。映像音声処理装置10において適用可能であるこのようなフィーチャは、非類似性測定を有していることが唯一の必要条件である。また、映像音声処理装置10は、実際には、フィーチャ抽出と上述したビデオ分割とを効率化のために同時に行うことがあり、以下に説明するフィーチャは、このような処理を可能にするものである。



## 【0038】

フィーチャとしては、まず映像フィーチャが挙げられる。セグメントを構成する映像（画像）は、そのセグメントの描写内容の大部分を表している。そのため、映像セグメントの類似性は、映像そのものの類似性に変わることができる場合が多い。したがって、映像フィーチャは、映像音声処理装置10で用いることができる重要なフィーチャの1種である。この映像フィーチャは、動的な情報よりも静的な情報を表すため、映像音声処理装置10は、後述するような方法によって、動的情報を得るように映像セグメント内の映像フィーチャを抽出する。

## 【0039】

映像フィーチャとして既知のものは、多数存在するが、シーン検出のためには以下に示す色フィーチャ（ヒストグラム）及び映像相関が、計算コストと精度との良好な兼ね合いを与えることを見出したことから、映像音声処理装置10は、映像フィーチャとしてこの色フィーチャ及び映像相関を用いる。

## 【0040】

映像音声処理装置10において、映像における色は、2つの映像が類似しているかを判断する際の重要な材料となる。カラーヒストグラムを用いて映像の類似性を判断することは、例えば“G. Ahanger and T.D.C. Little, A survey of technologies for parsing and indexing digital video, J. of Visual Communication and Image Representation 7:28-4, 1996”に記載されているように、よく知られている。カラーヒストグラムは、例えばHSVやRGB等の3次元色空間をn個の領域に分割し、各領域にある映像における画素の相対的割合を計算したものである。そして、得られた情報からは、n次元ベクトルが与えられる。圧縮されたビデオデータについては、例えばU.S. Patent #5,708,767号公報に記載されているように、カラーヒストグラムは、圧縮データから直接抽出される。

## 【0041】

映像音声処理装置10は、セグメントを構成する映像におけるもともとのYUV色空間を、色チャンネル当たり2ビットでサンプルし、長さ $2^{2 \cdot 3} = 64$ のヒストグラムを得る。

## 【0042】

このようなヒストグラムは、映像の全体的な色調を表すが、空間的情報を欠いている。そこで、映像音声処理装置10においては、映像相関を映像フィーチャの1つとして計算する。映像音声処理装置10におけるシーン検出に関しては、インターリーブされたセグメントは、シーンストラクチャの有力な指標を与える。例えば会話場面において、カメラの位置は、2人の話し手の間を交互に移動するが、カメラは、通常、同一の話し手を再度撮影するときには同じ位置に戻る。このような場合を検出するためには、映像のグレイスケールの相関がセグメントの類似性の良好な指標となることを見出したことから、映像音声処理装置10は、映像相関を計算するために、映像を $M \times N$ の大きさのグレイスケール映像にサブサンプルする。ここで、 $M$ と $N$ は、両方とも小さい値であり、例えば $8 \times 8$ である。小さいグレイスケール映像は、長さ $MN$ のフィーチャとして解釈される。

## 【0043】

さらに上述した映像フィーチャの他のフィーチャとしては、音声フィーチャが挙げられる。音声フィーチャは、音声セグメントの内容を表すのに用いることができるものである。映像音声処理装置10は、音声フィーチャとして周波数解析、ピッチ、レベルを用いる。これらの音声フィーチャは、種々の文献により知られているものである。

## 【0044】

まず、映像音声処理装置10は、フーリエ変換等の周波数解析を行うことによって、単一の音声フレームにおける周波数情報の分布を決定することができる。映像音声処理装置10は、例えば、1つの音声セグメントにわたる周波数情報の分布を表すために、FFT (Fast Fourier Transform; 高速フーリエ変換) 成分、周波数ヒストグラム、パワースペクトル、その他のフィーチャを用いることができる。

## 【0045】

また、映像音声処理装置10は、平均ピッチや最大ピッチといったピッチや、平均ラウドネスや最大ラウドネスといったレベルを音声セグメントを表す有効な音声フィーチャとして用いることができる。

## 【0046】

さらに他のフィーチャとしては、映像音声共通フィーチャが挙げられる。映像音声共通フィーチャは、特に映像フィーチャでもなく音声フィーチャでもないが、映像音声処理装置10において、シーン内のセグメントの内容を表すのに有用な情報を与えるものである。映像音声処理装置10は、この映像音声共通フィーチャとして、セグメント長とアクティビティとを用いる。

## 【0047】

映像音声処理装置10は、映像音声共通フィーチャとして、セグメント長を用いる。このセグメント長は、セグメントにおける時間長である。映像音声処理装置10において、シーンは、よくその変化のリズムの特性を有し、シーン内のセグメント長の傾向と同様の傾向をとる。例えば、迅速に連なった短いセグメントは、コマーシャルを表す。ところが、会話場面におけるセグメントは、より長くなるが、その長さの中でお互いに類似している。映像音声処理装置10は、このような特性を有するセグメント長を映像音声共通フィーチャとして用いることができる。

## 【0048】

また、映像音声処理装置10は、映像音声共通フィーチャとして、アクティビティを用いる。アクティビティは、セグメントの内容がどの程度動的或いは静的であるように感じられるかを表す。例えば、アクティビティは、視覚的に動的である場合には、カメラが対象物に沿って迅速に移動する度合若しくは撮影されている物が迅速に変化する度合を表す。

## 【0049】

このアクティビティは、カラーヒストグラムのようなフィーチャの平均フレーム間非類似性を測定することにより間接的に計算される。ここで、フレーム*i*とフレーム*j*との間で測定されたフィーチャ*F*に対する非類似性測定基準を $d_F(i, j)$ と定義すると、映像アクティビティ $V_F$ は、次式(1)のように定義される。

【0050】

【数1】

$$V_F = \frac{\sum_{i=b}^{f-1} d_F(i, i+1)}{f-b+1} \quad \dots (1)$$

【0051】

この式(1)において、bとfは、それぞれ、1セグメントにおける最初と最後のフレームのフレーム番号である。映像音声処理装置10は、具体的には、例えば上述したヒストグラムを用いて、映像アクティビティ $V_F$ を計算する。

【0052】

ところで、上述した映像フィーチャを始めとするフィーチャは、基本的にセグメントの静的情報を表すものであることは上述した通りであるが、セグメントの特徴を正確に表すためには、動的情報を表す必要がある。そこで、映像音声処理装置10は、以下に示すようなフィーチャのサンプリング方法により動的情報を表す。

【0053】

映像音声処理装置10は、例えば図5に示すように、1セグメント内の異なる時点から1以上の静的なフィーチャ値の組を抽出する。このとき、映像音声処理装置10は、フィーチャの抽出数を、忠実度の最大化と冗長度の最小化との釣り合いをとることにより決定する。例えば、セグメント内の1画像がキーフレームとして指定された場合には、そのキーフレームについて計算されたヒストグラムが抽出されたフィーチャとなる。

【0054】

映像音声処理装置10は、サンプリング方法を用いて、セグメントにおいて、あるフィーチャがとる全ての値のうち、どのサンプルを選択するかを決定する。映像音声処理装置10は、優れたサンプリング方法を必要とする。

【0055】

ここで、あるサンプルが常に所定の時点、例えばセグメント内の最後の時点に

おいてとられる場合を考える。この場合、黒にフェードする任意の2つのセグメントについては、サンプルが同一の黒いフレームとなるため、同一のフィーチャ値が得られる結果になる恐れがある。すなわち、これらのセグメントの映像内容がいかなるものであれ、サンプルした2つのフレームは、極めて類似していると判断されてしまう。このような問題は、サンプルが良好な代表値でないために発生するものである。

#### 【0056】

そこで、映像音声処理装置10は、このように固定点でフィーチャを抽出するのではなく、統計的な代表値を抽出する。ここでは、一般的なフィーチャのサンプリング方法を2つの場合、すなわち、第1の場合として、フィーチャを実数の $n$ 次元ベクトルとして表すことができる場合と、第2の場合として、非類似性測定基準しか利用できない場合とについて説明する。なお、第1の場合は、ヒストグラムやパワースペクトル等、最もよく知られている映像フィーチャ及び音声フィーチャを含むものである。

#### 【0057】

第1の場合には、サンプル数は、事前に $k$ と決められており、映像音声処理装置10は、“L. Kaufman and P.J. Rousseeuw, Finding Groups in Data: An Introduction to Cluster Analysis, John-Wiley and sons, 1990”に記載されてよく知られている $k$ 平均値クラスタリング法 ( $k$ -means-clustering method) を用いて、セグメント全体についてのフィーチャ値を $k$ 個の異なるグループに自動的に分割する。そして、映像音声処理装置10は、 $k$ 個の各グループから1サンプルを選択する。すなわち、映像音声処理装置10は、グループのセントロイド (平均ベクトル) 又はグループのセントロイドに近いサンプルを選択する。映像音声処理装置10は、この処理を短時間で行うことができ、サンプル数において直線的な時間しか要さない。

#### 【0058】

一方、第2の場合には、映像音声処理装置10は、“L. Kaufman and P.J. Rousseeuw, Finding Groups in Data: An Introduction to Cluster Analysis, John-Wiley and sons, 1990”に記載されている $k$ -メドイドアルゴリズム法 ( $k$ -me

doids algorithm method) を用いて、 $k$  個のグループを形成する。そして、映像音声処理装置 10 は、各グループ毎に、サンプル値として、上述したグループのセントロイドに類似するグループのメドイドを用いる。

#### 【0059】

なお、映像音声処理装置 10 においては、抽出されたフィーチャについての非類似性測定基準を構成する方法は、その基礎となるフィーチャの非類似性測定基準に基づくが、これについては後述する。

#### 【0060】

このようにして、映像音声処理装置 10 は、静的なフィーチャを抽出し、動的情報を表すことができる。

#### 【0061】

以上のように、映像音声処理装置 10 は、種々のフィーチャを抽出することができる。これらの各フィーチャは、一般に、単一ではセグメントの内容を表すのに不十分であることが多い。そこで、映像音声処理装置 10 は、これらの各種フィーチャを組み合わせることで、互いに補完し合うフィーチャの組を選択することができる。例えば、映像音声処理装置 10 は、上述したカラーヒストグラムと映像相関とを組み合わせることによって、各フィーチャが有する情報よりも多くの情報を得ることができる。

#### 【0062】

つぎに、図 4 中ステップ S3 におけるフィーチャを用いたセグメントの類似性測定について説明する。映像音声処理装置 10 は、2 つのフィーチャについて、その 2 つのフィーチャ値がどの程度非類似であるかを測定する実数値を計算する関数である非類似性測定基準を用いて、フィーチャ類似性測定部 17 によりセグメントの類似性測定を行う。この非類似性測定基準は、その値が小さい場合は 2 つのフィーチャが類似していることを示し、値が大きい場合は非類似であることを示す。ここでは、フィーチャ  $F$  に関する 2 つのセグメント  $S_1$ ,  $S_2$  の非類似性を計算する関数を非類似性測定基準  $d_F(S_1, S_2)$  と定義する。このような関数は、以下の式 (2) に示す特性に合致するものでなければならない。

【0063】

【数2】

$$\begin{aligned}
 d_F(S_1, S_2) &= 0 && (S_1 = S_2 \text{ のとき}) \\
 d_F(S_1, S_2) &\geq 0 && (\text{全ての } S_1, S_2 \text{ について}) \\
 d_F(S_1, S_2) &= d_F(S_2, S_1) && (\text{全ての } S_1, S_2 \text{ について})
 \end{aligned}
 \quad \dots (2)$$

【0064】

なお、適切な非類似性測定基準が、特定のフィーチャに依存することがあるが、“G. Ahanger and T.D.C. Little, A survey of technologies for parsing and indexing digital video, J. of Visual Communication and Image Representation 7:28-4, 1996” や “L. Kaufman and P.J. Rousseeuw, Finding Groups in Data: An Introduction to Cluster Analysis, John-Wiley and sons, 1990” に記載されているように、多くの一般的な非類似性測定基準は、 $n$ 次元空間における点として表されるフィーチャについての類似性を測定するのに有効であるとして知られている。これらは、ユークリッド距離、内積、 $L1$ 距離等である。特に、 $L1$ 距離が、ヒストグラムや映像相関といったフィーチャを含む種々のフィーチャに関して、有効に作用することを見出したことから、映像音声処理装置10は、2つの $n$ 次元ベクトル $A$ 、 $B$ 間の $L1$ 距離を次式(3)のように求める。

【0065】

【数3】

$$d_{L1}(A, B) = \sum_{i=1}^n |A_i - B_i| \quad \dots (3)$$

【0066】

ここで、下付文字は、 $n$ 次元ベクトルの $i$ 番目の要素を示すものである。

【0067】

また、映像音声処理装置10は、上述したように、時間とともに変化するフィ

ーチャに関して、セグメントにおける様々な時点でフィーチャ値を抽出する。そして、映像音声処理装置 10 は、2つの抽出されたフィーチャの間の類似性を決定するために、抽出されたフィーチャの非類似性を、その基礎となるフィーチャの非類似性測定基準により定義する。多くの場合、非類似性は、2つの抽出されたフィーチャのそれぞれから選択された最も非類似性のないフィーチャの対の非類似性値を用いて決定されるのが最良である。すなわち、映像音声処理装置 10 は、最小非類似性を決定する。この場合、2つの抽出されたフィーチャ  $SF_1$ 、 $SF_2$  の間のサンプルされた非類似性測定基準は、次式 (4) のように定義される。

【0068】

【数4】

$$d_s(SF_1, SF_2) = \min_{F_1 \in SF_1, F_2 \in SF_2} d_F(F_1, F_2) \quad \dots (4)$$

【0069】

上式 (4) における関数  $d_F(F_1, F_2)$  は、その基礎となる抽出されたフィーチャ  $F$  についての非類似性測定基準を示す。また、非類似性の値の最小をとる代わりに、最大又は平均をとる方がより適切である場合もある。

【0070】

さらに、映像音声処理装置 10 は、多くの場合、類似性を決定するために、同一セグメントについての多数のフィーチャから情報を組み合わせる必要がある。この1つの方法として、映像音声処理装置 10 は、種々のフィーチャベクトル非類似性関数の重み付きの組み合わせを計算する。すなわち、映像音声処理装置 10 は、 $k$  個のフィーチャ  $F_1, F_2, \dots, F_k$  が存在する場合、次式 (5) に表される組み合わせたフィーチャに関する非類似性測定基準  $d_F(S_1, S_2)$  を用いる。



【0071】

【数5】

$$d_F(S_1, S_2) = \sum_{i=1}^t w_i d_{Fi}(S_1, S_2) \quad \dots (5)$$

【0072】

ここで、 $\{w_i\}$  は、重みと  $\sum_i w_i = 1$  との組である。

【0073】

以上のように、映像音声処理装置10は、図4中ステップS2において抽出したフィーチャを用いて非類似性測定基準を計算し、セグメントの類似性を測定することができる。

【0074】

つぎに、図4中ステップS4におけるセグメントのグループ化について説明する。映像音声処理装置10は、非類似性測定基準と抽出したフィーチャとを用いて、時間的に近接して類似したセグメントを繰り返しまとめ、これらのセグメントを重ねてグループに合わせ、最終的に生成されたグループを検出シーンとして出力する。

【0075】

映像音声処理装置10は、セグメントをグループ化してシーンを検出する際に、2つの基本的な処理を行う。映像音声処理装置10は、まず第1の処理として、互いに時間的に近接して類似したセグメントのグループを検出する。この処理により得られるグループは、ほとんどが同一シーンの一部となるものである。また、映像音声処理装置10は、セグメントが時間的に重複しているため、第2の処理として、互いに時間が重複して類似したシーンをまとめる。映像音声処理装置10は、このような処理を各セグメントから開始し、反復して繰り返す。そして、映像音声処理装置10は、徐々にセグメントのグループを大きく構築していき、最終的に生成したグループをシーンの組として出力する。

## 【0076】

このような処理において、映像音声処理装置10は、その処理動作を制御するために2つの制約を用いる。

## 【0077】

すなわち、映像音声処理装置10は、第1の制約として、2つのセグメントが、どの程度類似している場合に同一のシーンのものであるのに十分類似しているかと考えられるかを決定する非類似性閾値 $\delta_{sim}$ を用いる。例えば、図6に示すように、映像音声処理装置10は、あるセグメントに対して一方のセグメントが類似性領域に属するか非類似性領域に属するかを判断する。

## 【0078】

なお、映像音声処理装置10は、非類似性閾値 $\delta_{sim}$ をユーザにより設定するようにしてもよく、また、後述するように、自動的に決定してもよい。

## 【0079】

また、映像音声処理装置10は、第2の制約として、2つのセグメントが、どの程度広く分割されている場合に同一のシーンのものであるかとなお考えられ得るかを決定する時間閾値Tを用いる。例えば、図7に示すように、映像音声処理装置10は、時間閾値Tの範囲内で互いに近接して続いている類似した2つのセグメントA、Bを同一シーンにまとめるが、時間的に大きく離れていて時間閾値Tの範囲外である2つのセグメントB、Cをまとめることはない。このように、映像音声処理装置10は、この時間閾値Tによる時間制約があるために、互いに類似しているものの大きく離れているセグメントを同一シーンにまとめてしまうエラーを発生することがない。

## 【0080】

なお、この時間閾値Tとしては、6～8の場合が概して良好な結果を与えることを見出したことから、映像音声処理装置10は、基本的に、時間閾値Tを6～8として用いる。

## 【0081】

映像音声処理装置10は、類似セグメントのグループを求めるために、ここでは、“L. Kaufman and P.J. Rousseeuw, Finding Groups in Data: An Introduct

ion to Cluster Analysis, John-Wiley and sons, 1990”に記載されている階層的クラスタ分割方法 (hierarchical clustering method) を適合させて用いることにする。このアルゴリズムにおいては、2つのクラスタ  $C_1$ ,  $C_2$  間の非類似性測定基準  $d_C(C_1, C_2)$  について、次式 (6) に示すように、含まれる2つの要素間の最小非類似性として定義する。

【0082】

【数6】

$$d_C(C_1, C_2) \equiv \min_{S_1 \in C_1, S_2 \in C_2} dist_S(S_1, S_2) \quad \dots (6)$$

【0083】

なお、映像音声処理装置10においては、上式 (6) で示される最小関数を最大又は平均関数に容易に置換することができる。

【0084】

まず、映像音声処理装置10は、図8に示すように、ステップS11において、変数Nをセグメントの数に初期化する。この変数Nは、常に現在検出されているグループの数を示すものである。

【0085】

次に、映像音声処理装置10は、ステップS12において、クラスタの組を生成する。映像音声処理装置10は、初期時には、N個の各セグメントを異なるクラスタとする。そのため、初期時には、N個のクラスタが存在することになる。各クラスタは、 $C^{\text{start}}$  と  $C^{\text{end}}$  によって表されるその開始時と終了時とを示す特性を有する。クラスタの組の要素は、 $C^{\text{start}}$  から順番に整頓される。

【0086】

さらに、映像音声処理装置10は、ステップS13において、変数tを1に初期化し、ステップS14において、変数tが時間閾値Tよりも大きいかな否かを判別する。ここで、映像音声処理装置10は、変数tが時間閾値Tよりも大きい場合には、ステップS23へと処理を移行し、変数tが時間閾値Tよりも小さい場

合には、ステップ S 15 へと処理を移行する。ただし、ここでは、変数  $t$  が 1 であるため、映像音声処理装置 10 は、ステップ S 15 へと処理を移行する。

#### 【0087】

映像音声処理装置 10 は、ステップ S 15 において、非類似性測定基準  $d_c$  を計算し、 $N$  個のクラスタの中から最も類似した 2 つのクラスタを検出する。ただし、ここでは、変数  $t$  が 1 であるため、映像音声処理装置 10 は、隣接したクラスタ間の非類似性測定基準  $d_c$  を計算し、その中から最も類似したクラスタの対を検出する。

#### 【0088】

このような最も類似した 2 つのクラスタを検出する方法としては、対象となる全てのクラスタの対を走査することが考えられる。映像音声処理装置 10 は、時間閾値  $T$  による時間制約を有するため、対象とするクラスタの対の数を限定することができる。そのため、映像音声処理装置 10 は、 $t$  個のセグメントにより分割されているクラスタのみを走査すればよい。クラスタの組は、時間順に整頓されているため、映像音声処理装置 10 は、あるクラスタに対してその前方及び後方へと走査していき、 $t$  個のセグメントよりも離れたクラスタを走査すれば、このクラスタ以降の全てのセグメントを対象外として走査処理を終了することができる。

#### 【0089】

ここで、検出された 2 つのクラスタをそれぞれ  $C_i$ 、 $C_j$  と定義し、これらのクラスタ  $C_i$ 、 $C_j$  の間の非類似性の値を  $d_{ij}$  と定義する。

#### 【0090】

映像音声処理装置 10 は、ステップ S 16 において、非類似性値  $d_{ij}$  が非類似性閾値  $\delta_{sim}$  よりも大きいかな否かを判別する。ここで、映像音声処理装置 10 は、非類似性値  $d_{ij}$  が非類似性閾値  $\delta_{sim}$  よりも大きい場合には、ステップ S 21 へと処理を移行し、非類似性値  $d_{ij}$  が非類似性閾値  $\delta_{sim}$  よりも小さい場合には、ステップ S 17 へと処理を移行する。ここでは、非類似性値  $d_{ij}$  が非類似性閾値  $\delta_{sim}$  よりも小さいものとする。

## 【0091】

映像音声処理装置10は、ステップS17において、クラスタ $C_j$ をクラスタ $C_i$ に合わせる。すなわち、映像音声処理装置10は、クラスタ $C_j$ の要素の全てをクラスタ $C_i$ に加える。

## 【0092】

また、映像音声処理装置10は、ステップS18において、クラスタ $C_j$ をクラスタの組から除去する。なお、2つのクラスタ $C_i$ 、 $C_j$ を合わせることにより開始時 $C_i^{\text{start}}$ の値が変化した場合には、映像音声処理装置10は、開始時による順序を維持するために、クラスタの組の要素を再び並べ替える。

## 【0093】

さらに、映像音声処理装置10は、ステップS19において、変数Nから1を減じる。

## 【0094】

そして、映像音声処理装置10は、ステップS20において、変数Nが1であるか否かを判別する。ここで、映像音声処理装置10は、変数Nが1である場合には、ステップS23へと処理を移行し、変数Nが1でない場合には、ステップS15へと処理を移行する。ここでは、変数Nが1でないものとする。

## 【0095】

すると、映像音声処理装置10は、ステップS15において、再び非類似性測定基準 $d_c$ を計算し、 $N-1$ 個のクラスタの中から最も類似した2つのクラスタを検出する。ここでも、変数 $t$ が1であるため、映像音声処理装置10は、隣接したクラスタ間の非類似性測定基準 $d_c$ を計算し、その中から最も類似したクラスタの対を検出する。

## 【0096】

また、映像音声処理装置10は、ステップS16において、非類似性値 $d_{ij}$ が非類似性閾値 $\delta_{\text{sim}}$ よりも大きいかな否かを判別する。ここでも、非類似性値 $d_{ij}$ が非類似性閾値 $\delta_{\text{sim}}$ よりも小さいものとする。

## 【0097】

そして、映像音声処理装置10は、ステップS17乃至ステップS20の処理

を行う。

【0098】

映像音声処理装置10は、このような処理を繰り返し、変数Nが減算されていった結果、ステップS20において、変数Nが1であると判別した場合には、ステップS23において、単一のセグメントのみを含むクラスタを合わせる。すなわち、映像音声処理装置10は、この場合は、全てのセグメントが1つのクラスタにまとめられているため、この処理を行う必要はなく、一連の処理を終了する。

【0099】

さて、映像音声処理装置10は、ステップS16において、非類似性値 $d_{ij}$ が非類似性閾値 $\delta_{sim}$ よりも大きいと判別した場合には、ステップS21へと処理を移行するが、この場合には、ステップS21において、時間的に重複しているクラスタを繰り返し合わせる。すなわち、 $C_i$ の時間間隔 $[C_i^{start}, C_i^{end}]$ が、 $C_j$ の時間間隔 $[C_j^{start}, C_j^{end}]$ と相交している場合には、2つのクラスタ $C_i$ と $C_j$ は、重複しているため、映像音声処理装置10は、クラスタをその組の開始時により整頓すると、重複しているクラスタを検出し、1つに合わせることができる。

【0100】

そして、映像音声処理装置10は、ステップS22において、変数 $t$ に1を加算して $t=2$ とし、ステップS14へと処理を移行して変数 $t$ が時間閾値 $T$ よりも大きいかな否かを判別する。ここでも、変数 $t$ が時間閾値 $T$ よりも小さいものとし、映像音声処理装置10は、ステップS15へと処理を移行するものとする。

【0101】

映像音声処理装置10は、ステップS15において、非類似性測定基準 $d_C$ を計算し、現在存在する複数のクラスタの中から最も類似した2つのクラスタを検出する。ただし、ここでは、変数 $t$ が2であるため、映像音声処理装置10は、1つおきに離れているクラスタ間の非類似性測定基準 $d_C$ を計算し、その中から最も類似したクラスタの対を検出する。

## 【0102】

そして、映像音声処理装置10は、ステップS16において、1つおきに離れているクラスタ $C_i$ 、 $C_j$ の非類似性値 $d_{ij}$ が非類似性閾値 $\delta_{sim}$ よりも大きいかなを判別する。ここでも、非類似性値 $d_{ij}$ が非類似性閾値 $\delta_{sim}$ よりも大きいものとし、映像音声処理装置10は、ステップS21及びステップS22の処理を経て、変数 $t$ に1を加算して $t=3$ として再びステップS14以降の処理へと移行する。勿論、映像音声処理装置10は、変数 $t$ が3のときには、ステップS15において、2つおきに離れているクラスタ間の非類似性測定基準 $d_c$ を計算し、その中から最も類似したクラスタの対を検出する。

## 【0103】

映像音声処理装置10は、このような処理を繰り返し、変数 $t$ が加算されていた結果、ステップS14において、変数 $t$ が時間閾値 $T$ よりも大きいと判別すると、ステップS23へと処理を移行し、単一のセグメントのみを含むクラスタを合わせる。すなわち、映像音声処理装置10は、孤立しているクラスタを単一の要素のみを含むクラスタとみなし、このような一連のクラスタが存在している場合には、これらのクラスタを単一のクラスタに合わせる。この工程は、単一のシーンと類似性関連を有さないセグメントをまとめるものである。なお、映像音声処理装置10は、必ずしもこの工程を行う必要はない。

## 【0104】

そして、映像音声処理装置10は、ステップS22の処理を経て一連の処理を終了する。

## 【0105】

このような一連の処理によって、映像音声処理装置10は、複数のクラスタをまとめていき、検出シーンを生成することができる。

## 【0106】

なお、映像音声処理装置10は、非類似性閾値 $\delta_{sim}$ をユーザにより設定するようにしてもよく、自動的に決定してもよいことは上述した通りである。非類似性閾値 $\delta_{sim}$ は、固定された値を用いる場合には、その最適な値がビデオデータの内容に依存する。非類似性閾値 $\delta_{sim}$ は、例えば、変化に富んだ映像内容を有

する番組の場合には、高い値に設定される必要があり、変化が少ない映像内容を有する番組の場合には、より低い値に設定される必要がある。そして、非類似性閾値  $\delta_{sim}$  が高すぎる場合には、検出されるシーンが少なくなりすぎ、非類似性閾値  $\delta_{sim}$  が低すぎる場合には、検出されるシーンが多くなりすぎることになる。

## 【0107】

このようなことから、映像音声処理装置 10 は、最適な非類似性閾値  $\delta_{sim}$  を決定することが、その性能を左右する上で重要となる。そのため、映像音声処理装置 10 は、非類似性閾値  $\delta_{sim}$  をユーザにより設定する場合には、上述したことを考慮した上で設定する必要がある。また、映像音声処理装置 10 は、以下に示す方法により、有効な非類似性閾値  $\delta_{sim}$  を自動的に決定することもできる。

## 【0108】

すなわち、映像音声処理装置 10 は、1つの方法として、平均値や最頻値といった統計的なフィーチャを用いて、 $(n)(n-1)/2$  個のセグメント対間の類似性の違いを表す距離の分布から非類似性閾値  $\delta_{sim}$  を得る。例えば、映像音声処理装置 10 は、全てのセグメント対における類似性の違いを表す距離の平均値  $\mu$  とその標準偏差  $\sigma$  とを用いて、非類似性閾値  $\delta_{sim}$  を  $a\mu + b\sigma$  と設定する。ここで、 $a$  及び  $b$  は、ある固定定数であるが、この定数  $a$  及び  $b$  は、それぞれ、0.5 及び 0.1 に設定することが良好な結果を与えることを見出している。

## 【0109】

実際には、映像音声処理装置 10 は、全てのセグメント対について、それらの間の類似性の違いを表す距離を決定する必要はない。映像音声処理装置 10 は、全てのセグメント対を対象とする代わりに、平均値  $\mu$  及び標準偏差  $\sigma$  の真値を判断するための距離値の部分集合をランダムに抽出する。すなわち、映像音声処理装置 10 は、距離を計算する際に、2つのセグメントをランダムに選択する。映像音声処理装置 10 は、2つのセグメントを十分に選択した場合には、平均値  $\mu$  及び標準偏差  $\sigma$  が真値に非常に近い結果を得ることができ、適切な非類似性閾値  $\delta_{sim}$  を自動的に得ることができる。このような場合、映像音声処理装置 10 は、例えば、セグメントの全数を  $n$ 、任意の小さい定数を  $C$  として、 $Cn$  と表され



るような数のセグメント対を抽出することによって、適切な非類似性閾値  $\delta_{sim}$  を自動的に決定することができる。

#### 【0110】

また、映像音声処理装置 10 は、シーンを検出する際に、セグメントが同一グループであるかを決定するために、単一の非類似性測定基準をもちいるばかりではなく、重み付け関数を用いて、異種のフィーチャに関する多様な非類似性測定基準を組み合わせることができることは上述した通りである。映像音声処理装置 10 において、このようなフィーチャの重み付けは、試行錯誤の末得られるものであり、各フィーチャが質的に異なるタイプのものである場合には、通常、適切な重み付けを行うことは困難であるが、映像音声処理装置 10 は、例えば、カラーヒストグラムとテクスチャフィーチャとを組み合わせる場合、各フィーチャに関してシーンを検出し、検出された各シーンストラクチャを単一のシーンストラクチャに組み合わせることができる。ここで、各フィーチャに関してシーンを検出した結果をシーン層と称することにする。例えば、フィーチャとしてカラーヒストグラムとセグメント長とを用いる場合、映像音声処理装置 10 は、これらのフィーチャに対してシーンを検出し、カラーヒストグラムについてのシーン層と、セグメント長についてのシーン層とを得て、これらのシーン層を単一のシーンストラクチャに組み合わせることができる。

#### 【0111】

さらに、一般に、映像領域と音声領域とからの情報を組み合わせることはできないが、映像音声処理装置 10 は、質的に異なるタイプのフィーチャを組み合わせる場合と同様な方法により、映像領域と音声領域とから得られるシーン層を単一のシーンストラクチャに組み合わせることができる。

#### 【0112】

このような処理について説明する。ここでは、それぞれが類似性の 1 つの基準を表す  $k$  個のフィーチャ  $F_1, F_2, \dots, F_k$  があるものとし、各フィーチャ  $F_i$  に対応して、非類似性測定基準  $d_F^i$  と、非類似性閾値  $\delta_{sim}^i$  と、時間閾値  $T^i$  とがあるものとする。映像音声処理装置 10 は、これらの各フィーチャ  $F_i$  に対する非類似性測定基準  $d_F^i$  と、非類似性閾値  $\delta_{sim}^i$  と、時間閾値  $T^i$  とを用いて

シーン層の組  $X_i = \{X_i^j\}$  を検出する。また、映像音声処理装置 10 は、映像情報と音声情報とに対して分割的にシーン層を検出し、映像情報と音声情報とに関する 2 つの独立したシーン層  $X_i = \{X_i^j\}$  ( $i = 1, 2$ ) を生成するものとする。

#### 【0 1 1 3】

映像音声処理装置 10 は、異なるシーン層を単一のシーンストラクチャに組み合わせるため、シーン境界の組み合わせ方を決定する必要がある。このシーン境界は、お互いにそろっている保証はないものである。ここで、各シーン層に関して、シーン境界を示す一連の時間で表される境界点  $t_{i1}, t_{i2}, \dots, t_{i|X_i|}$  があるものとする。まず、映像音声処理装置 10 は、種々のシーン層を単一のグループに組み合わせるために、最初にあるシーン層を整列に関する基礎とするために選択する。そして、映像音声処理装置 10 は、最終的に組み合わせる生成するシーンストラクチャにおけるシーン境界かどうかを各境界点  $t_{i1}, t_{i2}, \dots, t_{i|X_i|}$  に関して決定する。

#### 【0 1 1 4】

ここで、シーン層の組の番号である各  $i$  に関して、 $B_i(t)$  を、 $i$  番目のシーン層  $X_i$  において、ある時間  $t$  で「近く」にシーン境界があるかどうかを示すブール関数とする。この「近く」の意味は、変化し、例えば、映像情報と音声情報とをそろえている場合には、0.5 秒であったりする。

#### 【0 1 1 5】

映像音声処理装置 10 は、各境界点  $t_j = t_{ij}, j = 1, \dots, |X_j|$  に関して、 $l = 1, \dots, k$  のそれぞれについて、関数  $B_l(t_j)$  の結果を計算する。この結果は、各独立したシーン層に関して、シーン層  $X_i$  において時間  $t$  の近くにシーン境界があるかどうかを示している。そして、映像音声処理装置 10 は、決定関数として、組み合わせられた決定において時間  $t$  がシーン境界であるかどうかを決定するための  $B_i(t_j)$  の値を用いる。

#### 【0 1 1 6】

このような決定関数の 1 つの単純な例は、 $B_i(t_j)$  が定数  $m$  以上の 1 に等しい値の総数を計算することであり、時間の点  $t$  は、最後のシーンストラクチャに

おけるシーン境界を明示する。特に、 $m=1$ の場合、これは、「oring」境界点と同義のものである。 $m=k$ の場合、全ての境界点をそろえる要求と同義である。

#### 【0117】

このようにして、映像音声処理装置10は、異なるシーン層を単一のシーンストラクチャに組み合わせることができる。

#### 【0118】

以上説明してきたように、本発明の実施の形態として示す映像音声処理装置10は、シーンストラクチャを抽出するものであって、異なるタイプのビデオデータの内容に適用できるものである。映像音声処理装置10は、実際のビデオデータの内容において試験されてきており、テレビドラマや映画といった他のビデオデータジャンルからも同様に、シーンストラクチャを復活可能であることが論証済みである。

#### 【0119】

また、映像音声処理装置10は、完全に自動的であり、上述した非類似性閾値や時間閾値を設定するために、ユーザの介入を必要とせず、ビデオデータの内容の変化に応じて、適切な閾値を自動的に決定することができる。

#### 【0120】

さらに、映像音声処理装置10は、ユーザが事前にビデオデータの内容のストラクチャを知る必要はない。

#### 【0121】

さらにまた、映像音声処理装置10は、非常に単純であり計算上の効率もよいため、セットトップボックスやデジタルビデオレコーダ、ホームサーバ等の家庭用電子機器にも適用することができる。

#### 【0122】

また、映像音声処理装置10は、シーンを検出した結果、ビデオブラウジングのための新たな高レベルアクセスの基礎を得ることができる。そのため、映像音声処理装置10は、セグメントではなくシーンといった高レベルのビデオストラクチャを用いてビデオデータの内容を映像化することにより、内容に基づいたビ

デオデータへのアクセスが可能になる。例えば、映像音声処理装置 10 は、シーンを表示することにより、ユーザは、番組の要旨をすばやく知ることができ、興味のある部分を迅速に見つけることができる。

#### 【0123】

さらに、映像音声処理装置 10 は、シーン検出の結果、ビデオデータの概要又はダイジェストを自動的に作成するための基盤が得られる。一般に、ビデオデータからのランダムな断片を組み合わせるのとは異なり、一貫した概要を作成するには、ビデオデータを、再構成可能な意味を持つ成分に分解することができる必要がある。映像音声処理装置 10 により検出されたシーンは、そのようなダイジェストを作成するための正常な基盤である。

#### 【0124】

なお、本発明は、上述した実施の形態に限定されるものではなく、例えば、セグメント間の類似性測定のために用いるフィーチャ等は、上述したもの以外でもよいことは勿論であり、その他、本発明の趣旨を逸脱しない範囲で適宜変更が可能であることはいうまでもない。

#### 【0125】

##### 【発明の効果】

以上詳細に説明したように、本発明にかかる信号処理方法は、供給された信号の内容の意味構造を反映するパターンを検出して解析する信号処理方法であって、信号を構成する連続したフレームのひと続きから形成されるセグメントから、その特性を表す少なくとも 1 つ以上の特徴情報を抽出する特徴抽出工程と、特徴情報を用いて、特徴情報のそれぞれ毎に、セグメントの対の間の類似性を測定する測定基準を算出して、この測定基準によりセグメントの対の間の類似性を測定する類似性測定工程と、特徴情報と測定基準とを用いて、セグメントのうち、互いの時間的距離が所定の時間閾値以内であるとともに、互いの類似性が所定の非類似性閾値以上である 2 つのセグメントを検出し、信号の内容の意味構造を反映し、時間的に連続するセグメントからなるシーンにまとめるグループ化工程とを備える。

## 【0126】

したがって、本発明にかかる信号処理方法は、信号において類似したセグメントを検出してシーンにまとめることができ、セグメントよりも高いレベルのストラクチャを抽出することができる。

## 【0127】

また、本発明にかかる映像音声処理装置は、供給されたビデオ信号の内容の意味構造を反映する映像及び／又は音声のパターンを検出して解析する映像音声処理装置であって、ビデオ信号を構成する連続した映像及び／又は音声フレームのひと続きから形成される映像及び／又は音声セグメントから、その特性を表す少なくとも1つ以上の特徴情報を抽出する特徴抽出手段と、特徴情報を用いて、特徴情報のそれぞれ毎に、映像及び／又は音声セグメントの対の間の類似性を測定する測定基準を算出して、この測定基準により映像及び／又は音声セグメントの対の間の類似性を測定する類似性測定手段と、特徴情報と測定基準とを用いて、映像及び／又は音声セグメントのうち、互いの時間的距離が所定の時間閾値以内であるとともに、互いの類似性が所定の非類似性閾値以上である2つの映像及び／又は音声セグメントを検出し、ビデオ信号の内容の意味構造を反映し、時間的に連続する映像及び／又は音声セグメントからなるシーンにまとめるグループ化手段とを備える。

## 【0128】

したがって、本発明にかかる映像音声処理装置は、ビデオ信号において類似した映像及び／又は音声セグメントを検出してまとめ、シーンとして出力することが可能であり、映像及び／又は音声セグメントよりも高いレベルのビデオストラクチャを抽出することが可能となる。

## 【図面の簡単な説明】

## 【図1】

本発明において適用するビデオデータの構成を説明する図であって、モデル化したビデオデータのストラクチャを説明する図である。

## 【図2】

シーンを説明する図である。

【図 3】

本発明の実施の形態として示す映像音声処理装置の構成を説明するブロック図である。

【図 4】

同映像音声処理装置において、シーンを検出してグループ化する際の一連の工程を説明するフローチャートである。

【図 5】

同映像音声処理装置における動的フィーチャサンプリング処理を説明する図である。

【図 6】

非類似性閾値を説明する図である。

【図 7】

時間閾値を説明する図である。

【図 8】

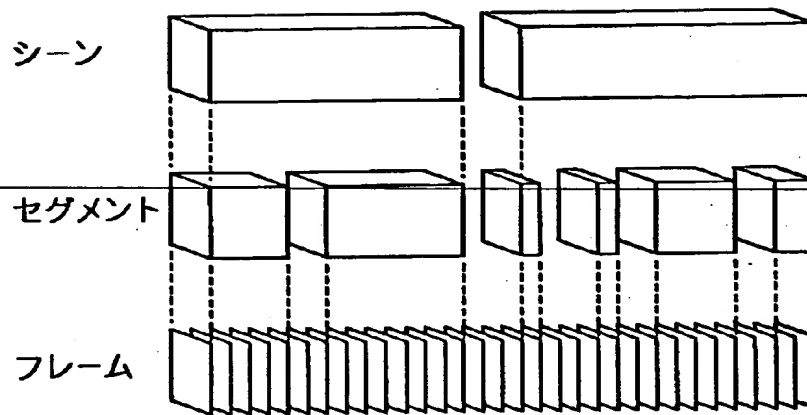
同映像音声処理装置において、セグメントをグループ化する際の一連の工程を説明するフローチャートである。

【符号の説明】

1 0 映像音声処理装置、 1 1 ビデオ分割部、 1 2 ビデオセグメントメモリ、 1 3 映像フィーチャ抽出部、 1 4 音声フィーチャ抽出部、 1 5 セグメントフィーチャメモリ、 1 6 シーン検出部、 1 7 フィーチャ類似性測定部

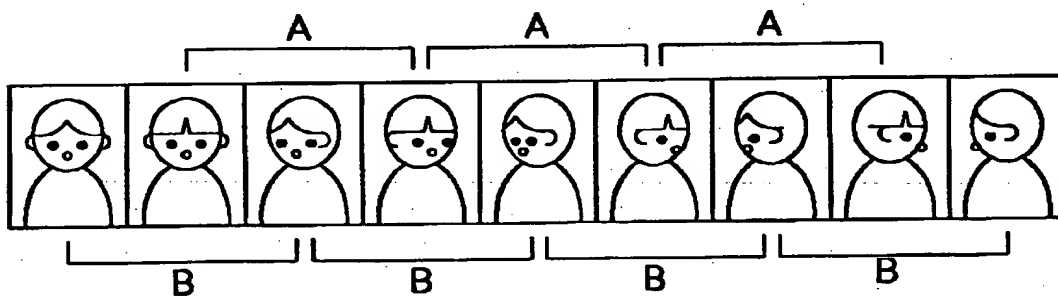
【書類名】 図面

【図 1】



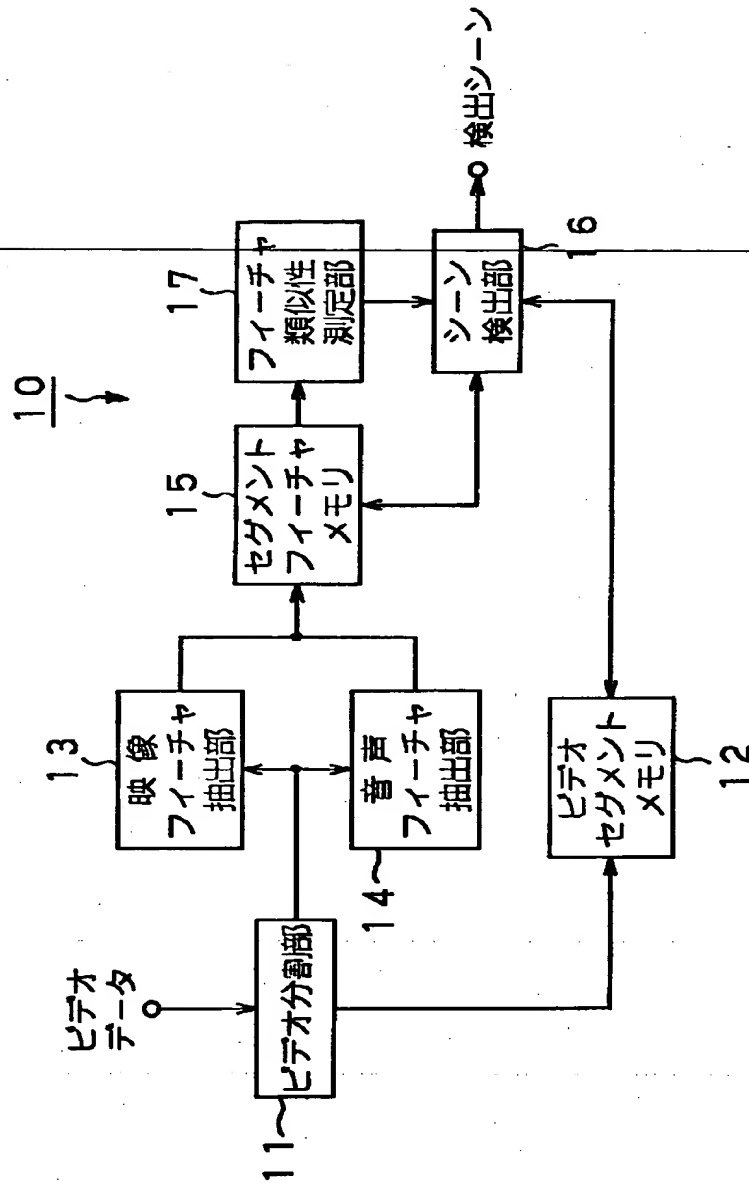
ビデオストラクチャの階層モデル

【図 2】



シーンの説明図

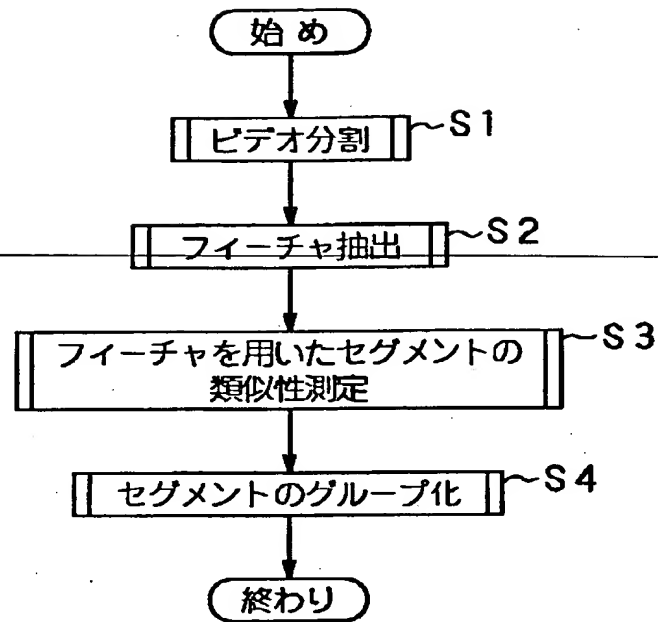
【図 3】



映像音声処理装置の構成ブロック図

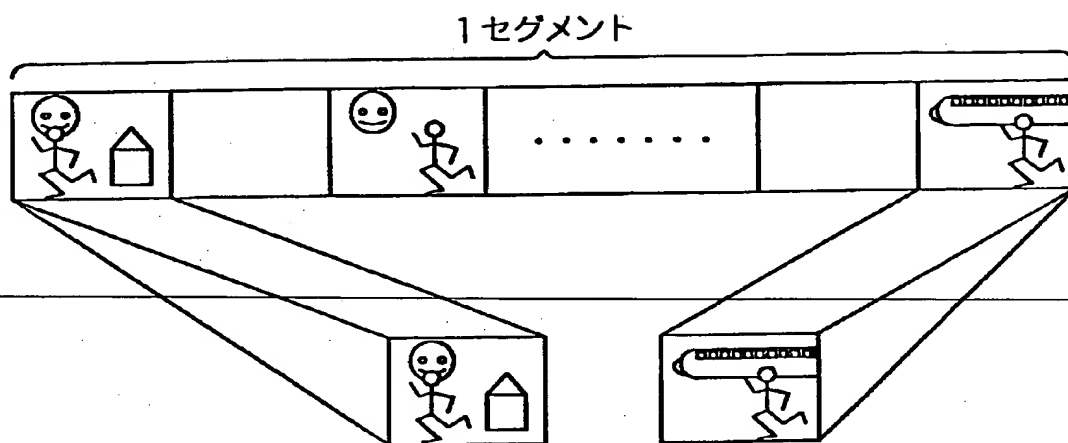


【図 4】



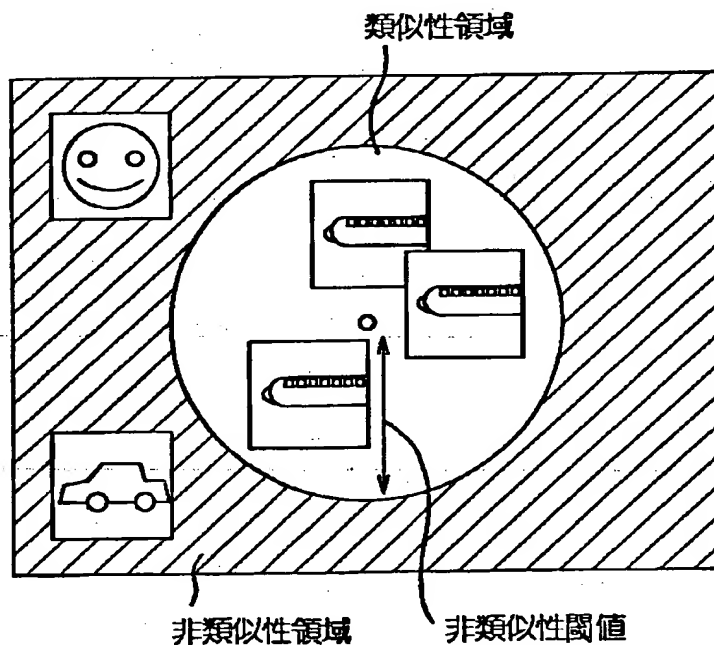
映像音声処理装置における一連の処理工程

【図 5】



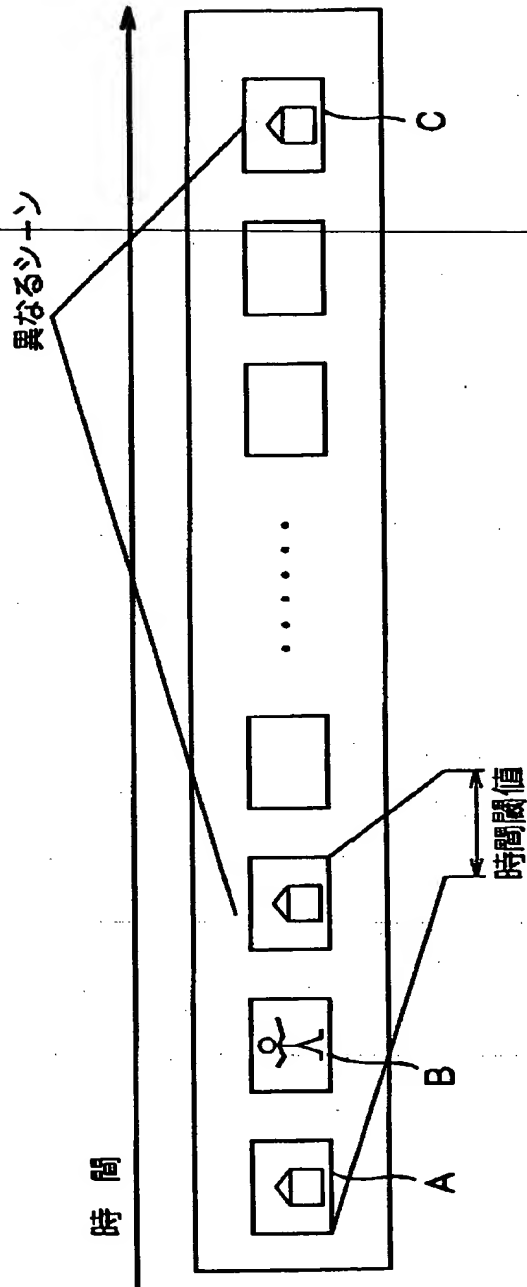
フィーチャのサンプリング方法の説明図

【図 6】



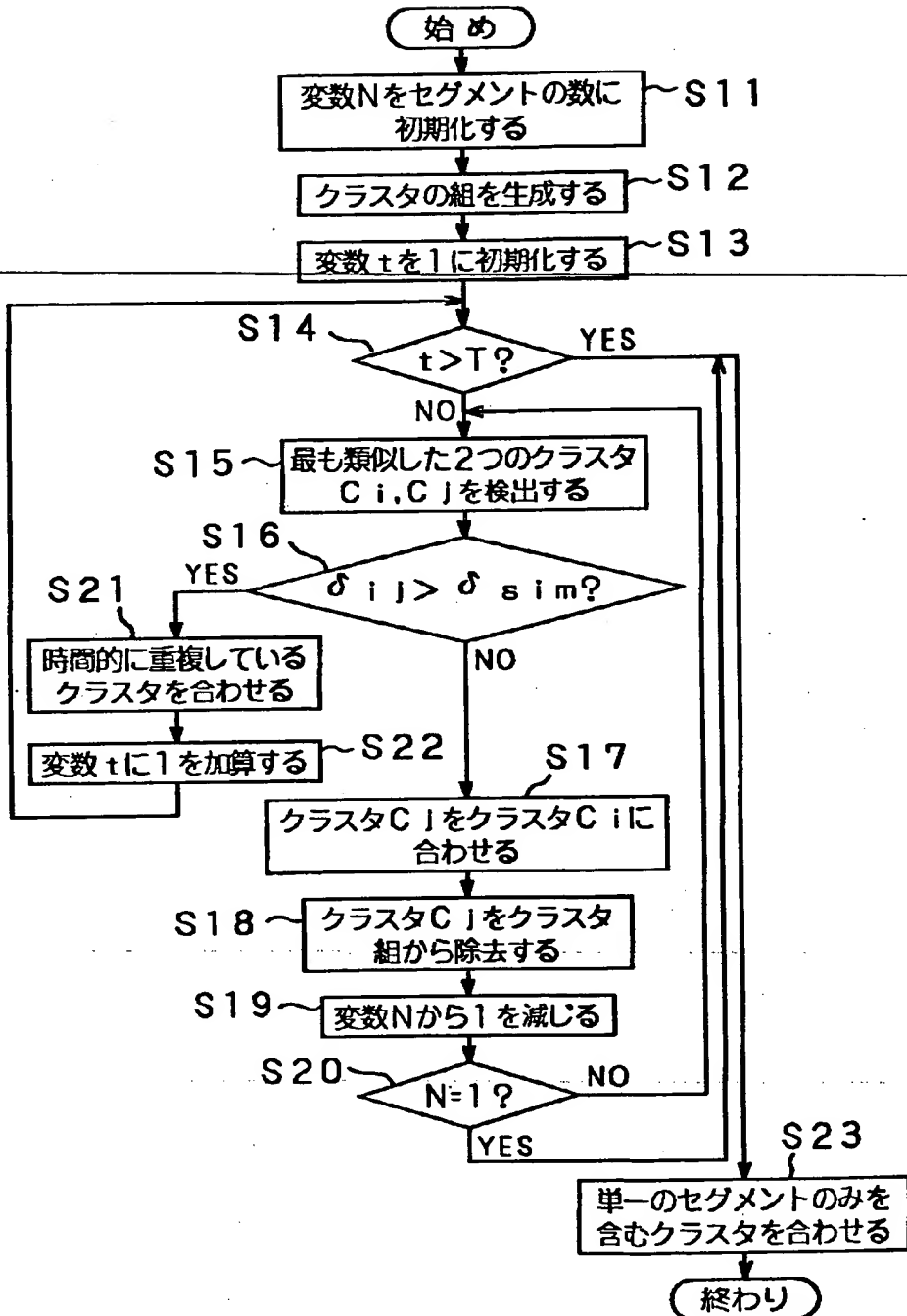
非類似性閾値の説明図

【図 7】



時間閾値の説明図

【図 8】



映像音声処理装置における一連の処理工程

【書類名】 要約書

【要約】

【課題】 種々のビデオにおける高レベルのビデオストラクチャを抽出する。

【解決手段】 映像音声処理装置 10 は、入力したビデオデータのストリームから分割された映像セグメント及び／又は音声セグメントから抽出されたフィーチャと、このフィーチャを用いて、各フィーチャ毎に計算された、映像セグメント及び／又は音声セグメントの対の間の類似性を測定する測定基準とを用いて、映像セグメント及び／又は音声セグメントのうち、互いの時間的距離が所定の時間閾値以内であるとともに、互いの類似性が所定の非類似性閾値以上である 2 つの映像セグメント及び／又は音声セグメントを検出し、ビデオデータの内容の意味構造を反映するシーンにまとめるシーン検出部 16 を備える。

【選択図】 図 3

特平11-023064

## 認定・付加情報

特許出願の番号	平成11年 特許願 第023064号
受付番号	59900079323
書類名	特許願
担当官	第一担当上席 0090
作成日	平成11年 2月13日

---

### <認定情報・付加情報>

【提出日】

平成11年 1月29日

次頁無

出 願 人 履 歴 情 報

識別番号

[000002185]

1. 変更年月日 1990年 8月30日

---

[変更理由] 新規登録

住 所 東京都品川区北品川6丁目7番35号

氏 名 ソニー株式会社

**THIS PAGE BLANK (USPTO)**